

# FSRNet: End-to-End Learning Face Super-Resolution with Facial Priors

Yu Chen<sup>1,5\*</sup> Ying Tai<sup>2\*</sup> Xiaoming Liu<sup>3</sup> Chunhua Shen<sup>4</sup> Jian Yang<sup>1</sup>  
<sup>1</sup>Nanjing University of Science and Technology <sup>2</sup>Youtu Lab, Tencent  
<sup>3</sup>Michigan State University <sup>4</sup>University of Adelaide <sup>5</sup>Motovis Inc.  
<https://github.com/tyshiwo/FSRNet>



Figure 1: Visual results of different super-resolution methods.

## Abstract

Face Super-Resolution (SR) is a domain-specific super-resolution problem. The facial prior knowledge can be leveraged to better super-resolve face images. We present a novel deep end-to-end trainable Face Super-Resolution Network (FSRNet), which makes use of the geometry prior, i.e., facial landmark heatmaps and parsing maps, to super-resolve very low-resolution (LR) face images without well-aligned requirement. Specifically, we first construct a coarse SR network to recover a coarse high-resolution (HR) image. Then, the coarse HR image is sent to two branches: a fine SR encoder and a prior information estimation network, which extracts the image features, and estimates landmark heatmaps/parsing maps respectively. Both image features and prior information are sent to a fine SR decoder to recover the HR image. To generate realistic faces, we also propose the Face Super-Resolution Generative Adversarial Network (FSRGAN) to incorporate the adversarial loss into FSRNet. Further, we introduce two related tasks, face alignment and parsing, as the new evaluation metrics for face SR, which address the inconsistency of classic metrics w.r.t. visual perception. Extensive experiments show that FSRNet and FSRGAN significantly outperforms state of the arts for very LR face SR, both quantitatively and qualitatively.

\*indicates equal contributions. This work was partially done when Yu Chen was visiting University of Adelaide. J. Yang is the corresponding author.

## 1. Introduction

Face Super-Resolution (SR), a.k.a. face hallucination, aims to generate a High-Resolution (HR) face image from a Low-Resolution (LR) input. It is a fundamental problem in face analysis, which can greatly facilitate face-related tasks, e.g., face alignment [16, 25, 36], face parsing [23], face recognition [34, 41], and 3D face reconstruction [29], since most existing techniques would degrade substantially when given very LR face images.

As a special case of general image SR, there exists face-specific prior knowledge in face images, which can be pivotal for face SR and is unavailable for general image SR [22, 32, 33]. For example, facial correspondence field could help recover accurate face shape [46], and facial components reveal rich facial details [31, 40]. However, as compared in Tab. 1, the previous face SR methods that utilize facial priors all adopt multi-stage, rather than end-to-end, training strategies, which is inconvenient and complicated.

Based on deep Convolutional Neural Network (CNN), in this work, we propose a novel *end-to-end trainable Face Super-Resolution Network (FSRNet)*, which estimates facial landmark heatmaps and parsing maps during training, and then uses these prior information to better super-resolve very LR face images. It is a consensus that end-to-end training is desirable for CNN [16], which has been validated in many areas, e.g., speech recognition [8] and image recognition [20]. Unlike previous Face SR methods that estimate local solutions in separate stages, our end-to-end framework learns the global solution directly, which is more convenient and elegant. To be specific, since it is non-trivial to estimate

Method	VDSR [17] (CVPR'16)	SRResNet [22] (CVPR'17)	StructuredFH [40] (CVPR'13)	CBN [46] (ECCV'16)	URDGN [42] (ECCV'16)	AttentionFH [2] (CVPR'17)	LCGE [31] (IJCAI'17)	FSRNet (ours)
Facial Prior KNWL	×	×	Components	Dense corres. field	×	×	Components	Landmark/parsing maps
Deep Model	✓	✓	×	✓	✓	✓	✓	✓
End-to-End	✓	✓	×	×	✓	✓	×	✓
Unaligned	✓	✓	×	✓	×	×	×	✓
Scale Factor	2/3/4	2/4	4	2/3/4	8	4/8	4	8

Table 1: **Comparisons with previous state-of-the-art super-resolution methods.** VDSR and SRResNet are generic image SR methods. StructuredFH, CBN, URDGN, AttentionFH and LCGE are face SR methods.

facial landmarks and parsing maps directly from LR inputs, we first construct a coarse SR network to recover a coarse HR image. Then, the coarse HR image is sent to a fine SR network, where a *fine SR encoder* and a *prior estimation network* share the coarse HR image as the input, followed by a *fine SR decoder*. The fine SR encoder extracts the image features, while the prior estimation network estimates landmark heatmaps and parsing maps jointly, via multi-task learning. After that, the image features and facial prior knowledge are fed into a fine SR decoder to recover the final HR face. The coarse and fine SR networks constitute our basic FSRNet, which already significantly outperforms the state of the arts (Fig. 1). To further generate realistic HR faces, *Face Super-Resolution Generative Adversarial Network (FSRGAN)* is introduced to incorporate the adversarial loss into the basic FSRNet. As in Fig. 1, FSRGAN recovers more realistic textures than FSRNet, and clearly shows superiority over the others.

It’s a consensus that Generative Adversarial Network (GAN)-based models recover visually plausible images but may suffer from low Peak Signal-to-Noise Ratio (PSNR), Structural SIMilarity (SSIM) or other quantitative metrics, while Mean Squared Error (MSE)-based deep models recover smooth images but with high PSNR/SSIM. To quantitatively show the superiority of GAN-based model, in [22], 26 users conducted a mean opinion score testing. However, such a testing is not objective and difficult to follow for fair comparison. To address this problem, we introduce two related face analysis tasks, *face alignment and parsing*, as the new evaluation metrics for face SR, which are demonstrated to be suitable for both MSE and GAN-based models.

In summary, the main contributions of this work include:

- To the best of our knowledge, this is the *first* deep face super-resolution network utilizing *facial geometry prior* in a convenient and elegant *end-to-end training* manner.
- Two kinds of facial geometry priors: *facial landmark heatmaps* and *parsing maps* are introduced simultaneously.
- The proposed FSRNet achieves the state of the art when hallucinating *unaligned* and *very low-resolution* ( $16 \times 16$  pixels) face images by an upscaling factor of 8, and the extended FSRGAN further generates more realistic faces.
- Face alignment and parsing are adopted as the *novel evaluation metrics* for face super-resolution, which are further demonstrated to resolve the inconsistency of classic metrics w.r.t. the visual perception.

## 2. Related Work

We review the prior works from two perspectives, and contrast with the most relevant papers in Tab. 1.

**Facial Prior Knowledge** There are many face SR methods that use facial prior knowledge to better super-resolve LR faces. Early techniques assume that faces are in a controlled setting with small variations [38]. Baker and Kanade [1] proposed to learn a prior on the spatial distribution of the image gradient for frontal face images. Wang et al. [37] implemented the mapping between LR and HR faces by an eigen transformation. Kolouri et al. [18] learnt a nonlinear Lagrangian model for HR face images, and enhanced the degraded image by finding the model parameters that could best fit the given LR data. Yang et al. [40] incorporated the face priors by using the mapping between specific facial components. However, the matchings between components are based on the landmark detection results that are difficult to estimate when the down-sampling factor is large.

Recently, deep convolutional neural networks have been successfully applied to the face SR task. Zhu et al. [46] super-resolved very LR and unaligned faces in a task-alternating cascaded framework. In their framework, face hallucination and dense correspondence field estimation are optimized alternatively. Besides, Song et al. [31] proposed a two-stage method, which first generated facial components by CNNs and then synthesized fine-grained facial structures through a component enhancement method. Different from the above methods that conduct face SR in multiple steps, our FSRNet fully leverages facial landmark heatmaps and parsing maps in an end-to-end training manner.

**End-to-end Training** End-to-end training is widely used in general image SR. Tai et al. [32] proposed Deep Recursive Residual Network (DRRN) to address the issue of model parameters and accuracy, which recursively learns the residual unit in a multi-path model. The authors also proposed a deep end-to-end persistent memory network to address the long-term dependency problem in CNN for image restoration [33]. Moreover, Ledig et al. [22] proposed Super-Resolution Generative Adversarial Network (SRGAN) for photo-realistic image SR using a perceptual loss function that consists of an adversarial loss and a content loss.

There are also many face SR methods adopting the end-to-end training strategy. Yu et al. [42] investigated GAN [7] to create perceptually realistic HR face images. The au-

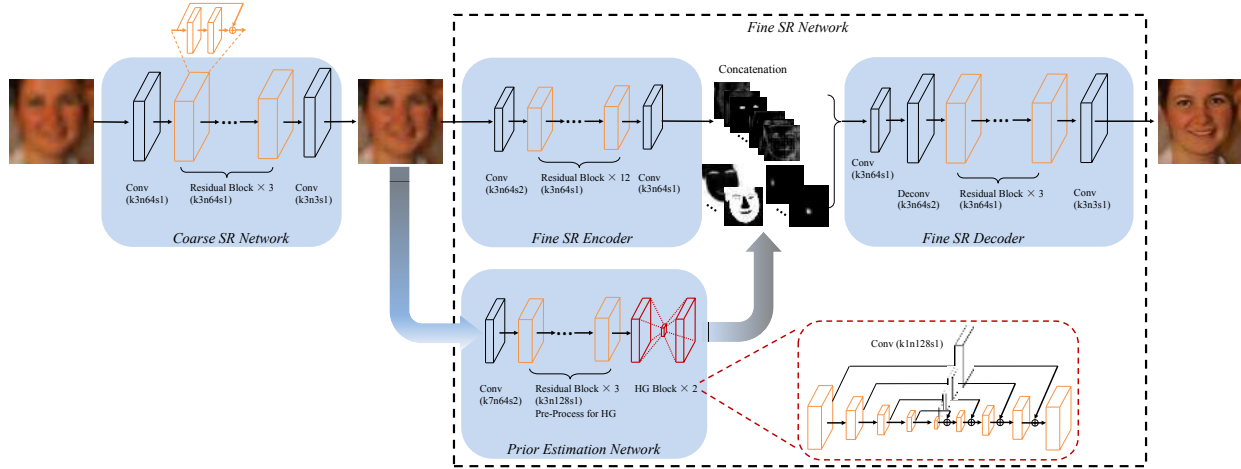


Figure 2: Network structure of the proposed FSRNet. ‘Conv’ indicates a convolutional layer with “pre-activation” structure [11]. ‘k3n64s1’ indicates the kernel size to be  $3 \times 3$ , the feature map number to be 64 and the stride to be 1.

thors further proposed transformative discriminative auto-encoder to super-resolve unaligned, noisy and tiny LR face images [43]. More recently, Cao et al. [2] proposed an attention-aware face hallucination framework, which resorts to deep reinforcement learning for sequentially discovering attended patches and then performing the facial part enhancement by fully exploiting the global image interdependency. Different from the above methods that only rely on the power of deep models, our FSRNet is not only an end-to-end trainable Neural Network, but also combines the rich information from the facial prior knowledge.

### 3. Face Super-Resolution Network

#### 3.1. Overview of FSRNet

Our basic FSRNet  $\mathbf{F}$  consists of four parts: *coarse SR network*, *fine SR encoder*, *prior estimation network* and finally a *fine SR decoder*. Denote  $\mathbf{x}$  as the low-resolution input image,  $\mathbf{y}$  and  $\mathbf{p}$  as the recovered high-resolution image and estimated prior information by FSRNet.

Since the very low-resolution input image may be too indistinct for prior estimation, we first construct the coarse SR network to recover a coarse SR image,

$$\mathbf{y}_c = \mathcal{C}(\mathbf{x}), \quad (1)$$

where  $\mathcal{C}$  denotes the mapping from a LR image  $\mathbf{x}$  to a coarse SR image  $\mathbf{y}_c$  by the coarse SR network. Then,  $\mathbf{y}_c$  is sent to the prior estimation network  $\mathcal{P}$  and fine SR encoder  $\mathcal{F}$ , as,

$$\mathbf{p} = \mathcal{P}(\mathbf{y}_c), \mathbf{f} = \mathcal{F}(\mathbf{y}_c), \quad (2)$$

where  $\mathbf{f}$  is the features extracted by  $\mathcal{F}$ . After encoding, the SR decoder  $\mathcal{D}$  is utilized to recover the SR image by *concatenating* the image feature  $\mathbf{f}$  and prior information  $\mathbf{p}$ ,

$$\mathbf{y} = \mathcal{D}(\mathbf{f}, \mathbf{p}). \quad (3)$$

Given a training set of  $N$  samples  $\{\mathbf{x}^{(i)}, \tilde{\mathbf{y}}^{(i)}, \tilde{\mathbf{p}}^{(i)}\}_{i=1}^N$ , where  $\tilde{\mathbf{y}}^{(i)}$  is the ground-truth HR image of the LR image  $\mathbf{x}^{(i)}$  and  $\tilde{\mathbf{p}}^{(i)}$  is the corresponding ground-truth prior information, FSRNet has the loss function,

$$\mathcal{L}_{\mathbf{F}}(\Theta) = \frac{1}{2N} \sum_{i=1}^N \{ \|\tilde{\mathbf{y}}^{(i)} - \mathbf{y}_c^{(i)}\|^2 + \alpha \|\tilde{\mathbf{y}}^{(i)} - \mathbf{y}^{(i)}\|^2 + \beta \|\tilde{\mathbf{p}}^{(i)} - \mathbf{p}^{(i)}\|^2 \}, \quad (4)$$

where  $\Theta$  denotes the parameter set,  $\alpha$  and  $\beta$  are the weights of the coarse SR loss and prior loss, and  $\mathbf{y}^{(i)}$ ,  $\mathbf{p}^{(i)}$  are the recovered HR image and estimated prior information of the  $i$ -th image respectively.

#### 3.2. Details inside FSRNet

We now present the details of our FSRNet, which consists of a coarse and a fine SR network, where the fine SR network contains three parts: a prior estimation network, a fine SR encoder and a fine SR decoder.

##### 3.2.1 Coarse SR network

First, we use a coarse SR network to roughly recover a coarse HR image. The motivation is that it is non-trivial to estimate facial landmark positions and parsing maps directly from a LR input image. Using the coarse SR network may help to ease the difficulties for estimating the priors. The architecture of the coarse SR network is shown in Fig. 2. It starts with a  $3 \times 3$  convolution followed by 3 *residual blocks* [10]. Then another  $3 \times 3$  convolutional layer is used to reconstruct the coarse HR image.

##### 3.2.2 Fine SR Network

In the following fine SR network, the coarse HR image is sent to two branches, prior estimation network and fine en-

coder network, to estimate facial priors and extract features, respectively. Then the decoder jointly uses results of both branches to recover the fine HR image.

**Prior Estimation Network** Any real-world object has distinct distributions in its shape and texture, including face. Comparing facial shape with texture, we choose to model and leverage the shape prior for two considerations. First, when reducing the resolution from high to low, the shape is better preserved compared to the texture, and hence is more likely to be extracted to facilitate super-resolution. Second, it is much easier to represent shape prior than texture prior. E.g., face parsing estimates the segmentations of different face components, and landmarks provide the accurate locations of facial keypoints, even at low resolution [26]. Both represent facial shapes, while parsing carries more granularity. In contrast, it is not clear how to represent the higher-dimensional texture prior for a specific face.

Inspired by the recent success of stacked heatmap regression in human pose estimation [3, 28], we adopt the Hour-Glass (HG) structure to estimate facial landmark heatmaps and parsing maps in our prior estimation network. Since both priors represent the 2D face *shape*, in our prior estimation network, *the features are all shared between these two tasks*, except the last layer. The detailed structure of prior estimation network is shown in Fig. 2. To effectively consolidate features across scales and preserve spatial information in different scales, the hourglass block uses a skip connection mechanism between symmetrical layers. An  $1 \times 1$  convolution layer follows to post-process the obtained features. Finally, the shared hourglass feature is connected to two separate  $1 \times 1$  convolution layers to generate the landmark heatmaps and the parsing maps.

**Fine SR Encoder** For fine SR encoder, inspired by the success of ResNet [10] in SR [22, 32], we utilize the residual blocks for feature extraction. Considering the computation cost, the size of our prior features is down-sampled to  $64 \times 64$ . To make the feature size consistent, the fine SR encoder starts with a  $3 \times 3$  convolutional layer of stride 2 to down-sample the feature map to  $64 \times 64$ . Then the ResNet structure is utilized to extract image features.

**Fine SR Decoder** The fine SR decoder jointly uses the features and priors to recover the final fine HR image. First, the prior feature  $\mathbf{p}$  and image feature  $\mathbf{f}$  are concatenated as the input of the decoder. Then a  $3 \times 3$  convolutional layer reduces the number of feature maps to 64. A  $4 \times 4$  deconvolutional layer is utilized to up-sample the feature map to size  $128 \times 128$ . Then 3 residual blocks are used to decode the features. Finally, a  $3 \times 3$  convolutional layer is used to recover the fine HR image.

### 3.3. FSRGAN

As we know, GAN has shown great power in super-resolution [22], which generates photo-realistic images with

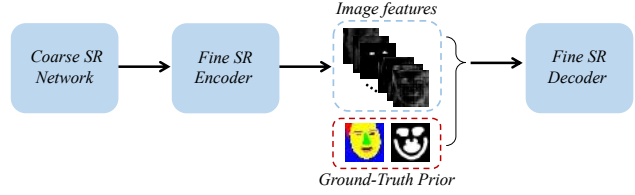


Figure 3: **Structure of “upper-bound” model.** The ground-truth priors are directly concatenated with image features. Removing priors in the red box and increasing the number of image features by the number of channels in prior induce to the baseline model.

superior visual effect than MSE-based deep models. The key idea is to use a discriminative network to distinguish the super-resolved images and the real high-resolution images, and to train the SR network to deceive the discriminator.

To generate realistic high-resolution faces, our model utilizes GAN in the conditional manner [13]. The objective function of the adversarial network  $\mathbf{C}$  is expressed as:

$$\mathcal{L}_{\mathbf{C}}(\mathbf{F}, \mathbf{C}) = \mathbb{E}[\log \mathbf{C}(\tilde{\mathbf{y}}, \mathbf{x})] + \mathbb{E}[\log(1 - \mathbf{C}(\mathbf{F}(\mathbf{x}), \mathbf{x}))], \quad (5)$$

where  $\mathbf{C}$  outputs the probability of the input been real and  $\mathbb{E}$  is the expectation of the probability distribution. Apart from the adversarial loss  $\mathcal{L}_{\mathbf{C}}$ , we further introduce a perceptual loss [15] using high-level feature maps (i.e., features from ‘relu5\_3’ layer) of the pre-trained VGG-16 network [30] to help assess perceptually relevant characteristics,

$$\mathcal{L}_{\mathbf{P}} = \|\phi(\mathbf{y}) - \phi(\tilde{\mathbf{y}})\|^2, \quad (6)$$

where  $\phi$  denotes the *fixed* pre-trained VGG model, and maps the images  $\mathbf{y}/\tilde{\mathbf{y}}$  to the feature space. In this way, the final objective function of FSRGAN is:

$$\arg \min_{\mathbf{F}} \max_{\mathbf{C}} \mathcal{L}_{\mathbf{F}}(\Theta) + \gamma_{\mathbf{C}} \mathcal{L}_{\mathbf{C}}(\mathbf{F}, \mathbf{C}) + \gamma_{\mathbf{P}} \mathcal{L}_{\mathbf{P}}, \quad (7)$$

where  $\gamma_{\mathbf{C}}$  and  $\gamma_{\mathbf{P}}$  are the weights of GAN and perceptual loss, respectively.

## 4. Prior Knowledge for Face Super-Resolution

In this section, we would like to answer two questions: (1) Is facial prior knowledge really useful for face super-resolution? (2) How much improvement does different facial prior knowledge bring? To answer these questions, we conduct several tests on the 2,330-image Helen dataset [21]. The last 50 images are used for testing and the others are for training. We perform data augmentation on the training images. Specifically, we rotate the original images by  $90^\circ$ ,  $180^\circ$ ,  $270^\circ$  and flip them horizontally. This results in 7 additional augmented images for each original one. Besides, each image in Helen dataset has a ground truth label of 194 landmarks and 11 parsing maps.

**Effects of Facial Prior Knowledge** First, we demonstrate that facial prior knowledge is *significant* for face super-resolution, even without any advanced processing steps.

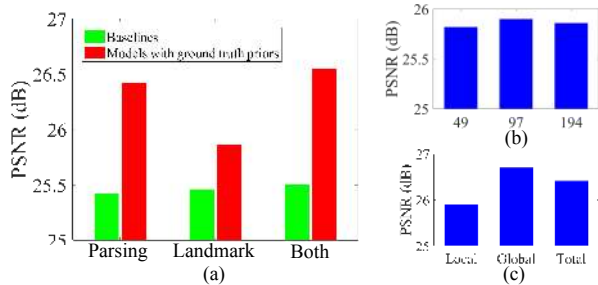


Figure 4: **Effects of facial priors.** (a) Comparison between baselines and models with ground truth priors. The upper bound performance of landmark priors with different numbers of landmarks (b), and parsing priors with different types of parsing maps (c).

We remove the prior estimation network and construct a single-branch baseline network. Based on the baseline network, we introduce the ground truth facial prior information (i.e., landmark heatmaps and parsing maps) to the “concatenation” layer to construct a new network, as shown in Fig. 3. For fair comparison, we keep the feature map number of “concatenation” layer the same between two networks, which means the results can contrast the effects of the facial prior knowledge. Fig. 4 presents the performance of 3 kinds of settings, including setting with or without parsing maps, landmark heatmaps, or both maps, respectively. As we can see, the models using prior information significantly outperform the corresponding baseline models with the PSNR improvement of 0.4 dB after using landmark heatmaps, 1.0 dB after using parsing maps, and 1.05 dB after using both priors, respectively. These huge improvements on PSNR clearly signify the *positive* effects of facial prior knowledge to face SR.

**Upper Bound Improvements from Priors** Next, we focus on specific prior information, and study the upper bound improvements that different priors bring. Specifically, for facial landmarks, we introduce 3 sets of landmarks, i.e., 49, 97 and 194 landmarks, respectively. For parsing maps, we introduce the global and local parsing maps, respectively. The global parsing map is shown in Figs. 5(b-c), while Fig. 5(d) shows the local parsing maps containing different facial components. From the results of different priors in Fig. 4, we observe that: (1) Parsing priors contain richer information for face SR and bring more improvements than the landmark prior. (2) Global parsing maps are more useful than local parsing maps. (3) More landmark heatmaps have minor improvements than the version using 49 landmarks.

The above results and analysis demonstrate the effects of both facial priors, and show the upper bound performance that we achieve if the priors are predicted *perfectly*. Since we use the recent popular facial alignment/parsing framework as the prior estimation network, the powerful learning ability enables the network to leverage the priors as much as possible, and hence can benefit the face SR. Apart from

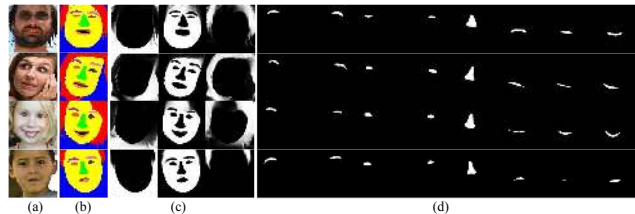


Figure 5: **Parsing maps of Helen images.** (a) Original image. (b) Color visualization map generated by 11 ground truth parsing maps [23]. It is used as part of the global parsing map. (c) Global parsing maps from the ground truth. (d) Local parsing maps from the ground truth, containing left eyebrow, right eyebrow, left eye, right eye, nose, upper lip, inner mouth, and lower lip, respectively.

the benefit to PSNR, introducing facial prior may bring other advantages, such as more precise recovery of the *face shape*, as reflected by less errors on face alignment and parsing. More details are presented in the next section.

## 5. Experiments

### 5.1. Implementation Details

**Datasets** We conduct extensive experiments on 2 datasets: Helen [21] and celebA [27]. Experimental setting on Helen dataset is described in Sec. 4. For celebA dataset, we use the first 18,000 images for training, and the following 100 images for evaluation. It should be noted that celebA only has a ground truth of 5 landmarks. We further use a recent alignment model [4] to estimate the 68 landmarks and adopt GFC [23] to estimate the parsing maps as the ground truth. **Training Setting** We coarsely crop the training images according to their face regions and resize to  $128 \times 128$  without any pre-alignment operation. For testing, any popular face detector [9] can be used to obtain the cropped image as the input. Same as [22], color images are used for training. The input low-resolution images are firstly enlarged by bicubic interpolation, and hence have the same size as the output high-resolution images. For implementation, we train our model with the Torch7 toolbox [5]. The model is trained using the RMSprop algorithm with an initial learning rate of  $2.5 \times 10^{-4}$ , and the mini-batch size of 14. We empirically set  $\alpha = 1$ ,  $\beta = 1$ ,  $\gamma_C = 10^{-3}$  and  $\gamma_P = 10^{-1}$  for both datasets. Training a basic FSRNet on Helen dataset takes  $\sim 6$  hours on 1 Titan X GPU.

### 5.2. Ablation Study

**Effects of Estimated Priors** We conduct ablation study on the effects of the prior estimation network. Since our SR branch has the similar network structure as SRResNet [22], we clearly show how the performance improves with different kinds of facial priors based on the performance of SRResNet. In this test, we *estimate* the facial priors through the prior estimation network instead of using the ground truth conducted in Sec. 4. Same as the tests conducted

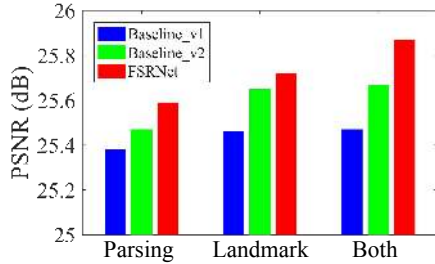


Figure 6: Ablation study on effects of estimated priors.

in Fig. 4 (a), we conduct 3 experiments to estimate the landmark heatmaps, parsing maps, or both maps, respectively. In each experiment, we further compare our basic FSRNet with two other network structures. Specifically, by removing the prior estimation network from our basic FSRNet, the remaining parts constitute the first network, named ‘Baseline\_v1’, which has the similar structure and hence similar performance as SRResNet. The second network, named ‘Baseline\_v2’, has the same structure as our basic FSRNet except that there is no supervision on the prior estimation network.

Fig. 6 shows the results of different network structures. It can be seen that: (1) The second networks always outperform the first networks. The reason may be even without the supervision, the second branch learns additional features that provide more high-frequency signals to help SR. (2) Compared to the second networks, the supervision on prior knowledge further improves the performance, which indicates the estimated facial priors indeed have positive effects on face super-resolution. (3) The model using both priors achieves the best performance, which indicates richer prior information brings more improvement. (4) The best performance reaches 25.87 dB, which is lower than the performance (i.e., 26.55 dB) when using ground truth. That means our estimated priors are not perfect and a better prior estimation network may result in higher model performance.

**Effects of Hourglass Numbers** As discussed in Sec. 4, a powerful prior estimation network may lead to accurate prior estimation. Here, we study the effect of the hourglass number  $h$  in the prior estimation network. Specifically, we test  $h = 1/2/4$ , and the PSNR results are 25.69, 25.87, and 25.95 dB, respectively. Since using more hourglasses leads to a deeper structure, the learning ability of the prior estimation network grows, and hence better performance. To intuitively show the adjustments in stacking more hourglasses, we show the landmark estimations of the first and second stacked hourglass in Fig 7. It can be observed that the estimation is obviously improved in the second stacking.

**Effects of End-to-end Training** Next, we show that end-to-end training helps both prior estimation and face SR. Specifically, we train coarse SR, prior branch and fine SR of FSRNet separately on Helen dataset, which achieve 24.21 dB, 5.61 NRMSE and 25.65 dB respectively (vs. 24.26 dB,



CFAN	CFSS	SDM	DeepAlign	FSRNet_S1	FSRNet_S2
9.45	7.26	7.88	6.50	9.44	7.04

Figure 7: Landmark estimations by FSRNet on CelebA. **First row:** Results of the first stacked HG (FSRNet\_S1). **Second row:** Results of the second HG (FSRNet\_S2). Please zoom in to see the improvements. **Bottom:** NRMSEs of the first four methods are achieved by testing directly on the ground-truth HR images.

5.28 NRMSE and 25.87 dB of FSRNet). End-to-end training obviously contributes to performance improvement.

### 5.3. Comparisons with State-of-the-Art Methods

We compare FSRNet with state-of-the-art SR methods, including generic SR methods like SRResNet [22], VDSR [17] and SRCNN [6]; and facial SR methods like GLN [35] and URDGN [42]. For fair comparison, we use the released codes of the above models and train all models with the same training set. For URDGN [42], we only train the generator to report PSNR/SSIMs, but the entire GAN network for qualitative comparisons.

**Face Super-Resolution** First, we compare FSRNet with the state of the arts quantitatively. Tab. 2 summarizes quantitative results on the two datasets. Our FSRNet significantly outperforms state of the arts in both PSNR and SSIM. Not surprisingly, FSRGAN achieves low PSNR/SSIMs. Besides, we also present FSRNet\_aug, which sends multiple augmented test images during inference and then fuse the outputs to report the results. This simple yet effective trick brings significant improvements.

Qualitative comparisons of FSRNet/FSRGAN with prior works are illustrated in Fig. 8. Benefiting from the facial prior knowledge, our method produces relatively sharper edges and shapes, while other methods may give more blurry results. Moreover, FSRGAN further recovers sharper facial textures than FSRNet.

We next compare FSRGAN with two recent face SR methods: Wavelet-SRNet [12] and CBN [46]. We follow the same experimental setting on handling occluded face as [12] and directly import the  $16 \times 16$  test examples from [12] for super-resolving  $128 \times 128$  HR images. As shown in Fig. 9, FSRGAN achieves relatively sharper shapes (e.g., nose in all cases) than the state of the arts.

**Face Alignment** Apart from evaluating PSNR/SSIM, we introduce face alignment as a novel evaluation metric for face super-resolution, since accurate face recovery

Dataset	Bicubic	SRCNN	VDSR	SRResNet	GLN	URDGN	FSRNet	FSRNet_Laug	FSRGAN
Helen	23.69/0.6592	23.97/0.6779	24.61/0.6980	25.30/0.7297	24.11/0.6922	24.22/0.6909	<b>25.87/0.7602</b>	<b>26.21/0.7720</b>	25.10/0.7234
celebA	23.75/0.6423	24.26/0.6634	24.83/0.6878	25.82/0.7369	24.55/0.6867	24.63/0.6851	<b>26.31/0.7522</b>	<b>26.60/0.7628</b>	25.20/0.7023

Table 2: **Benchmark super-resolution**, with PSNR/SSIMs for scale factor 8. Red/blue color indicate the best/second best performance.

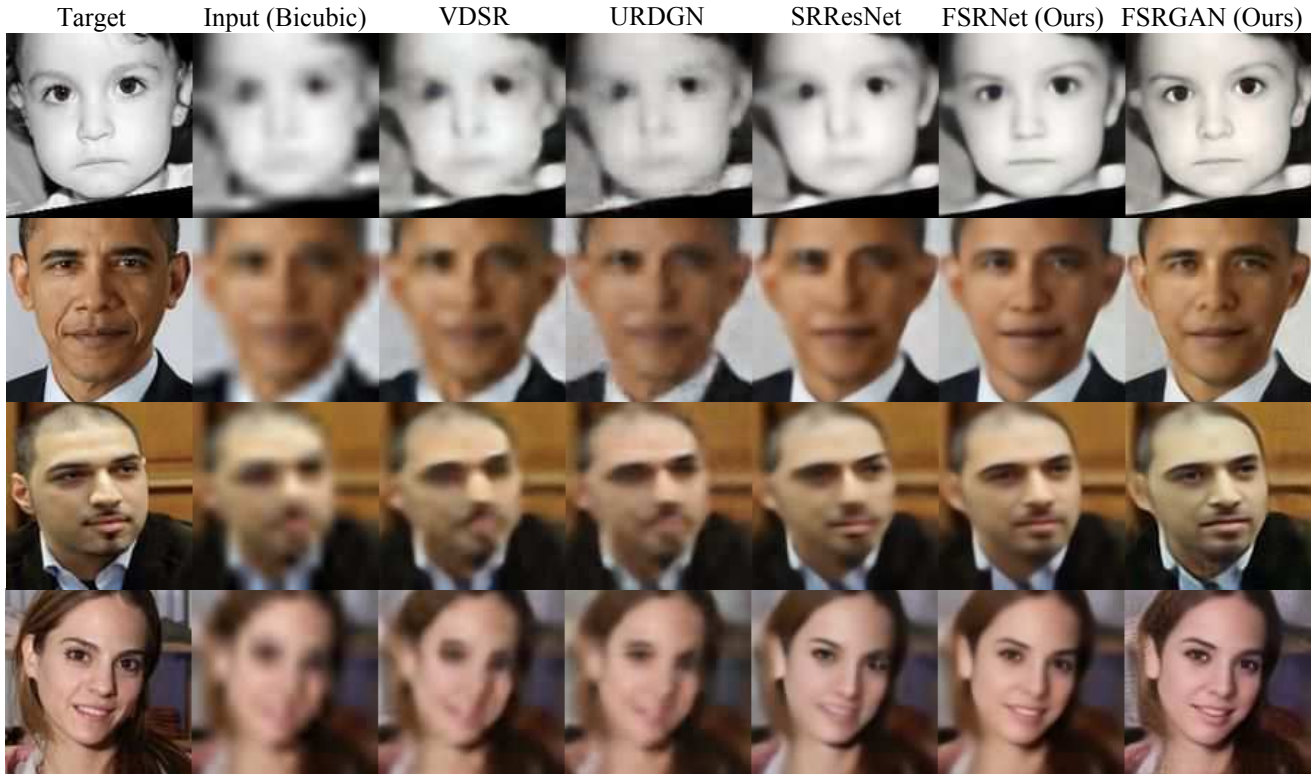


Figure 8: **Qualitative comparisons**. Top two examples are of Helen and others are of celebA. Please zoom in to see the differences.

Dataset	Bicubic	SRCNN	VDSR	SRResNet	GLN	URDGN	FSRNet	FSRNet_Laug	FSRGAN	Target
Helen	5.89/0.2908	5.58/0.3442	5.29/0.3691	4.87/0.4555	5.72/0.3694	5.22/0.4070	4.18/0.5758	<b>4.13/0.5817</b>	<b>3.94/0.6128</b>	3.32/0.6744
celebA	13.3/0.2319	12.7/0.2912	12.4/0.3329	11.3/0.5453	12.2/0.4058	12.2/0.3553	<b>10.6/0.6195</b>	<b>10.6/0.6269</b>	<b>10.2/0.6518</b>	9.45/—

Table 3: **Quantitative comparisons on alignment (NRMSE)/parsing (IoU)**. For celebA, parsing maps from target HR images are the GT.

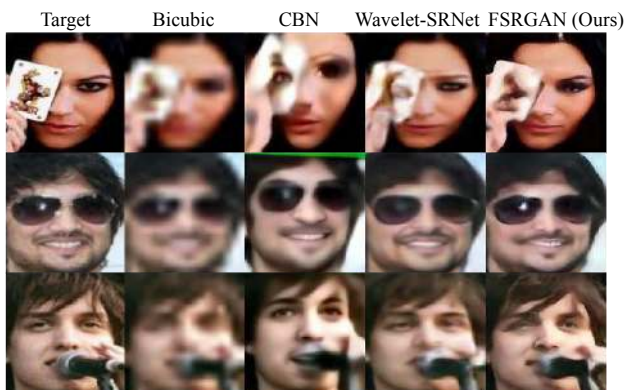


Figure 9: **Comparisons with CBN and Wavelet-SRNet**.

should lead to accurate shape/geometry, and hence accurate landmark points. We adopt a popular alignment model CFAN [44] to estimate the landmarks of different recovered images. The upper part of Fig. 10 shows the recovered im-

ages of SRResNet and our FSRNet, including the results from coarse SR net and final output. The bottom part shows the facial landmarks estimated by CFAN on different recovered images, which are directly displayed on the target image for clear comparisons. Tab. 3 also presents the Normalized Root Mean Squared Error (NRMSE) results, which is a popular metric in face alignment and lower NRMSE indicates better alignment performance. From the results we can see that: (1) It is challenge for the state-of-the-art alignment models to estimate landmarks directly from very low-resolution images. The estimated landmarks of the bicubic image exhibit large errors around mouth, eyes or other components. In FSRNet, the coarse SR net can ease the alignment difficulty to some extent, which leads to lower NRMSE than the input bicubic image. (2) Compared to SRResNet, our final output provides visually superior estimation on mouth, eyes and shape, and also achieves a large margin of 0.9 quantitatively. That demonstrates the effectiveness of using landmark priors for training.

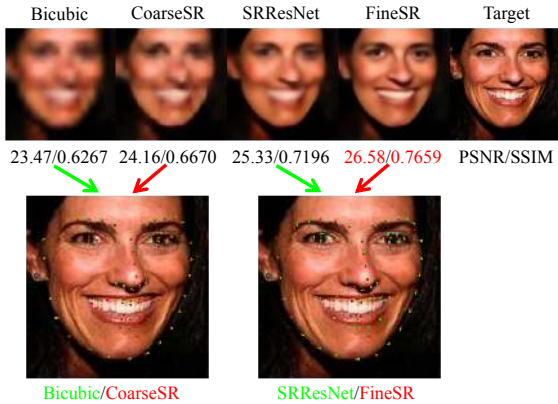


Figure 10: Qualitative comparison of face alignment.

**Face Parsing** We also introduce face parsing as another evaluation metric for face super-resolution. Although our prior estimation network can predict the parsing maps from the LR inputs, for fair comparison, we adopt a recent facial parsing model GFC [23] to generate the parsing maps for the recovered images of all methods, including the bicubic inputs, our coarse SR net, SRResNet, our fine SR net, and targets, respectively. Intersection-over-union (IoU) is reported in Tab. 3. As we can see, the coarse SR net also has positive effects on face parsing. Fig. 11 presents the estimated parsing maps by [23], the parsing maps from our final HR images recover complete and accurate components, while SRResNet may generate wrong shapes or even lose components (e.g., mouth).

Here, we adopt two side tasks, face alignment and parsing, as the new evaluation metrics for face super-resolution. They can subjectively evaluate the quality of geometry in the recovered images, which is complementary to the classic PSNR/SSIM metrics that focus more on photometric quality. Further, Tab. 3 shows that FSRGAN outperforms FSRNet on both metrics, which is consistent with the superior visual quality in Fig. 8. This consistency actually addresses one issue in GAN-based SR methods, which has the superior visual quality, but lower PSNR/SSIM. This also shows that GAN-based methods can better recover the facial geometry, in addition to perceived visual quality.

**Prior Estimation** Priors estimated by FSRNet are by-products of our model. Here, we first compare the landmarks *directly* estimated by FSRNet with methods [19, 39, 44, 45] using their released codes, as shown in the bottom of Fig. 7 and Fig. 12. It should be noted that *our method starts with the LR images while others are tested directly on the ground-truth  $8 \times HR$  images*. Despite the disadvantage in the input image resolution, our method outperforms most recent methods and is competitive with the state of the art [19]. It should also be noted that estimating 194 points is more difficult than 68 points [45], especially on the low-resolution faces. Then we also compare the IoU of parsing maps estimated by FSRNet with GFC in Fig. 12.

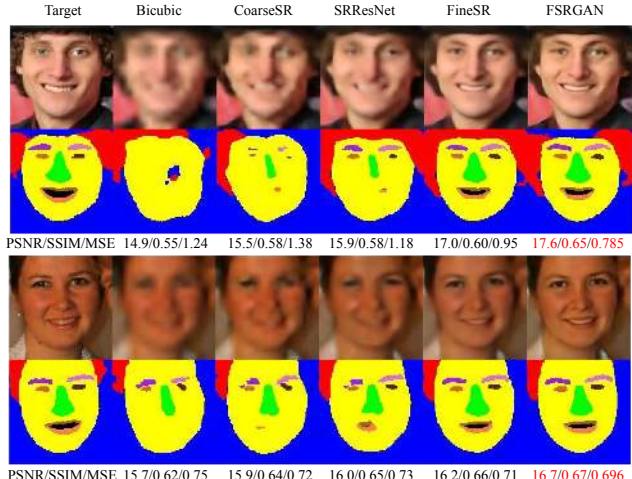


Figure 11: Qualitative comparison of face parsing.

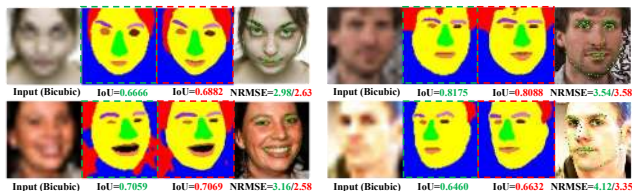


Figure 12: Visualizations of estimated prior. Green indicates results from prior branch. Red indicates results from final SR image, by GFC [23] for parsing and CFAN [44] for 68 landmarks.

**Time Complexity** Unlike CBN that needs multiple steps and trains multiple models for face hallucination, our FSRNet is lightweight, faster and more convenient to use, which only needs *one forward process* for inference and costs 0.012s on Titan X GPU, for a  $128 \times 128$  image. For comparison, CBN has four cascades and totally consumes 3.84s [46], while the traditional face SR requires more time, e.g., [24] needs 8 minutes and [14] needs 15 – 20 minutes.

## 6. Conclusions

In this paper, a novel deep end-to-end trainable Face Super-Resolution Network (FSRNet) is proposed for face super-resolution. The key component of FSRNet is the prior estimation network, which not only helps to improve the photometric recovery in terms of PSNR/SSIM, but also provides a solution for accurate geometry estimation directly from very LR images, as shown in the results of facial landmarks/parsing maps. Extensive experiments show that FSRNet is superior to the state of the arts on unaligned face images, both quantitatively and qualitatively. Following the main idea of this work, future research can be expanded in various aspects, including designing a better prior estimation network, e.g., learning the fine SR network iteratively, and investigating other useful facial priors, e.g., texture.



## References

- [1] S. Baker and T. Kanade. Hallucinating faces. In *Proc. IEEE Int. Conf. Automatic Face and Gesture Recognition.*, 2000. 2
- [2] Q. Cao, L. Lin, Y. Shi, X. Liang, and G. Li. Attention-aware face hallucination via deep reinforcement learning. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2017. 2, 3
- [3] Y. Chen, C. Shen, X.-S. Wei, L. Liu, and Y. Jian. Adversarial PoseNet: A Structure-aware Convolutional Network for Human Pose Estimation. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2017. 4
- [4] Y. Chen, C. Shen, X.-S. Wei, L. Liu, and J. Yang. Adversarial learning of structure-aware fully convolutional networks for landmark localization. *arXiv: Comp. Res. Repository*, 2017. 5
- [5] R. Collobert, K. Kavukcuoglu, and C. Farabet. Torch7: A matlab-like environment for machine learning. In *Proc. Advances in Neural Inf. Process. Syst., BigLearn Workshop*, 2011. 5
- [6] C. Dong, C. Loy, K. He, and X. Tang. Image super-resolution using deep convolutional networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 38(2):295–307, 2016. 6
- [7] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Proc. Advances in Neural Inf. Process. Syst.*, 2014. 2
- [8] A. Graves, A. rahman Mohamed, and G. Hinton. Speech recognition with deep recurrent neural networks. In *Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing.*, 2013. 1
- [9] Z. Hao, Y. Liu, H. Qin, J. Yan, X. Li, and X. Hu. Scale-aware face detection. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2017. 5
- [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016. 3, 4
- [11] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. In *Proc. Eur. Conf. Comp. Vis.*, 2016. 3
- [12] H. Huang, R. He, Z. Sun, and T. Tan. Wavelet-srnet: A wavelet-based CNN for multi-scale face super resolution. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2017. 6
- [13] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2017. 4
- [14] Y. Jin and C. Bouganis. Robust multi-image based blind face hallucination. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2015. 8
- [15] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *Proc. Eur. Conf. Comp. Vis.*, 2016. 4
- [16] A. Jourabloo, M. Ye, X. Liu, and L. Ren. Pose-invariant face alignment with a single CNN. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2017. 1
- [17] J. Kim, J. K. Lee, and K. M. Lee. Accurate image super-resolution using very deep convolutional networks. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2016. 2, 6
- [18] S. Kolouri and G. K. Rohde. Transport-based single frame super resolution of very low resolution face images. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2015. 2
- [19] M. Kowalski, J. Naruniec, and T. Trzcinski. Deep alignment network: A convolutional neural network for robust face alignment. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn., Faces-in-the-wild Workshop/Challenge*, 2017. 8
- [20] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Proc. Advances in Neural Inf. Process. Syst.*, 2012. 1
- [21] V. Le, J. Brandt, Z. Lin, L. Boudev, and T. S. Huang. Interactive facial feature localization. In *Proc. Eur. Conf. Comp. Vis.*, 2012. 4, 5
- [22] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-realistic single image super-resolution using a generative adversarial network. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2017. 1, 2, 4, 5, 6
- [23] Y. Li, S. Liu, J. Yang, and M.-H. Yang. Generative face completion. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2017. 1, 5, 8
- [24] C. Liu, H. Shum, and W. Freeman. Face hallucination: Theory and practice. *Int. J. Comput. Vision*, 2007. 8
- [25] F. Liu, D. Zeng, Q. Zhao, and X. Liu. Joint face alignment and 3D face reconstruction. In *Proc. Eur. Conf. Comp. Vis.*, 2016. 1
- [26] X. Liu, P. H. Tu, and F. W. Wheeler. Face model fitting on low resolution images. In *Proc. British Machine Vis. Conf.*, 2006. 4
- [27] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2015. 5
- [28] A. Newell, K. Yang, and J. Deng. Stacked hourglass networks for human pose estimation. In *Proc. Eur. Conf. Comp. Vis.*, 2016. 4
- [29] J. Roth, Y. Tong, and X. Liu. Adaptive 3D face reconstruction from unconstrained photo collections. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(11), November 2017. 1
- [30] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *Proc. Int. Conf. Learn. Representations*, 2015. 4
- [31] Y. Song, J. Zhang, S. He, L. Bao, and Q. Yang. Learning to hallucinate face images via component generation and enhancement. In *Proc. Int. Joint Conf. Artificial Intell.*, 2017. 1, 2
- [32] Y. Tai, J. Yang, and X. Liu. Image super-resolution via deep recursive residual network. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2017. 1, 2, 4
- [33] Y. Tai, J. Yang, X. Liu, and C. Xu. Memnet: A persistent memory network for image restoration. In *Proc. IEEE Int. Conf. Comp. Vis.*, 2017. 1, 2
- [34] Y. Taigman, M. Yang, M. Ranzato, and L. Wolf. Deepface: Closing the gap to human-level performance in face verification. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2014. 1
- [35] O. Tuzel, Y. Taguchi, and J. R. Hershey. Global-local face upsampling network. *arXiv: Comp. Res. Repository*, 2016. 6

- [36] G. Tzimiropoulos. Project-out cascaded regression with an application to face alignment. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2015. [1](#)
- [37] X. Wang and X. Tang. Hallucinating face by eigentransformation. *IEEE Trans. Syst., Man, Cybern. C*, 35(3):425–434, 2005. [2](#)
- [38] F. W. Wheeler, X. Liu, and P. H. Tu. Multi-frame super-resolution for face recognition. In *Proc. IEEE Int. Conf. Biometrics: Theory, Applications, and Systems*, 2007. [2](#)
- [39] X. Xiong and F. De la Torre. Supervised descent method and its applications to face alignment. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2013. [8](#)
- [40] C.-Y. Yang, S. Liu, and M.-H. Yang. Structured face hallucination. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2013. [1](#), [2](#)
- [41] J. Yang, L. Luo, J. Qian, Y. Tai, F. Zhang, and Y. Xu. Nuclear norm based matrix regression with applications to face recognition with occlusion and illumination changes. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(1):156–171, 2017. [1](#)
- [42] X. Yu and F. Porikli. Ultra-resolving face images by discriminative generative networks. In *Proc. Eur. Conf. Comp. Vis.*, 2016. [2](#), [6](#)
- [43] X. Yu and F. Porikli. Hallucinating very low-resolution unaligned and noisy face images by transformative discriminative autoencoders. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2017. [3](#)
- [44] J. Zhang, S. Shan, M. Kan, and X. Chen. Coarse-to-fine auto-encoder networks (CFAN) for real-time face alignment. In *Proc. Eur. Conf. Comp. Vis.*, 2014. [7](#), [8](#)
- [45] S. Zhu, C. Li, C. Change Loy, and X. Tang. Face alignment by coarse-to-fine shape searching. In *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2015. [8](#)
- [46] S. Zhu, S. Liu, C. Loy, and X. Tang. Deep cascaded bi-network for face hallucination. In *Proc. Eur. Conf. Comp. Vis.*, 2016. [1](#), [2](#), [6](#), [8](#)